Pearson's correlation coefficient

between

X and Y

for sample size n

definition

$$-1 \leq r = \frac{S_{XY}}{\sqrt{S_{XX}}\sqrt{S_{YY}}} \leq 1$$

Where

$$S_{XY} = \sum_{i=1}^{n} (Y_i - \bar{Y})(X_i - \bar{X})$$

$$= \sum_{i=1}^{n} X_i Y_i - n\bar{X}\bar{Y}$$

$$S_{XX} = \sum_{i=1}^{n} (X_i - \bar{X})^2$$

$$= \sum_{i=1}^{n} X_i^2 - n(\bar{X})^2$$

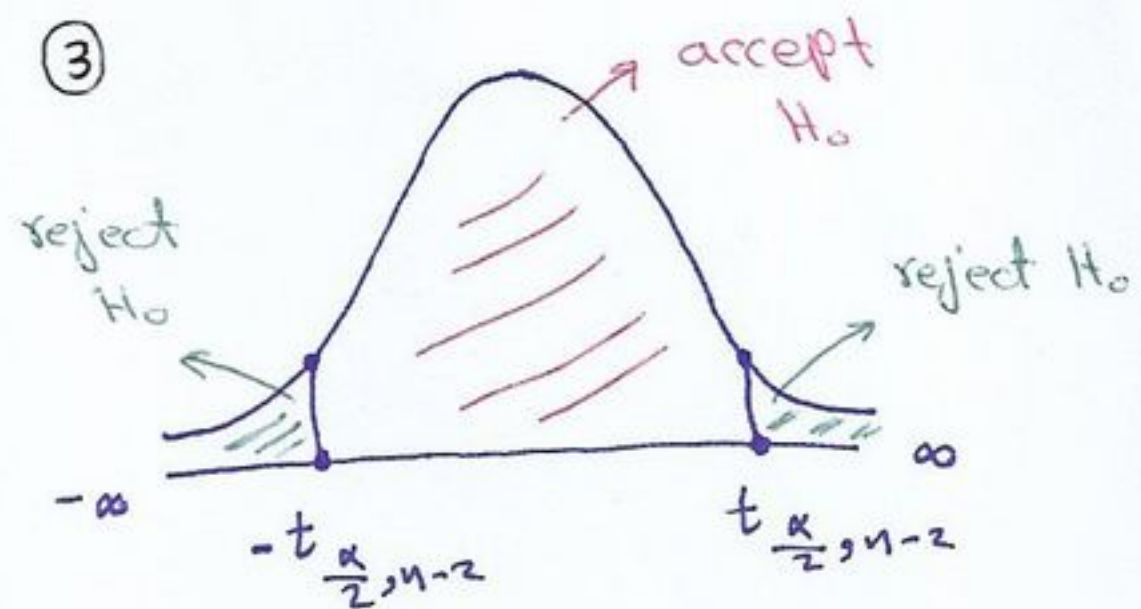$$S_{YY} = \sum_{i=1}^{n} (Y_i - \bar{Y})^2$$

$$= \sum_{i=1}^{n} Y_i^2 - n(\bar{Y})^2$$

Hypotheses testing

① $H_0 : \rho = 0$ Vs $H_1 : \rho \neq 0$

② test statistics

$$t_{T.s.} = \frac{r\sqrt{n-2}}{\sqrt{1-r^2}} \sim t_{v=n-2}$$

③



accept $H_0$

reject $H_0$

reject $H_0$

$-\infty$

$-t_{\frac{\alpha}{2}, n-2}$

$t_{\frac{\alpha}{2}, n-2}$

$\infty$

④ or → we reject $H_0$ if:

$$t_{T.s.} < -t_{\frac{\alpha}{2}, n-2}$$

or

$$t_{T.s.} > t_{\frac{\alpha}{2}, n-2}$$

otherwise, we will accept $H_0$

we reject $H_0$ if:

$$P\text{-value} \leq \alpha$$

where

$$P\text{-value} = 2 P(T > |t_{T.s.}|)$$

①

# Some notes:

① Pearson's correlation coefficient
   for study
   the straight line/linear relationship
   between
   X and Y

② $\rho \longmapsto$ correlat. for population

   $r \longmapsto$ correlation for sample
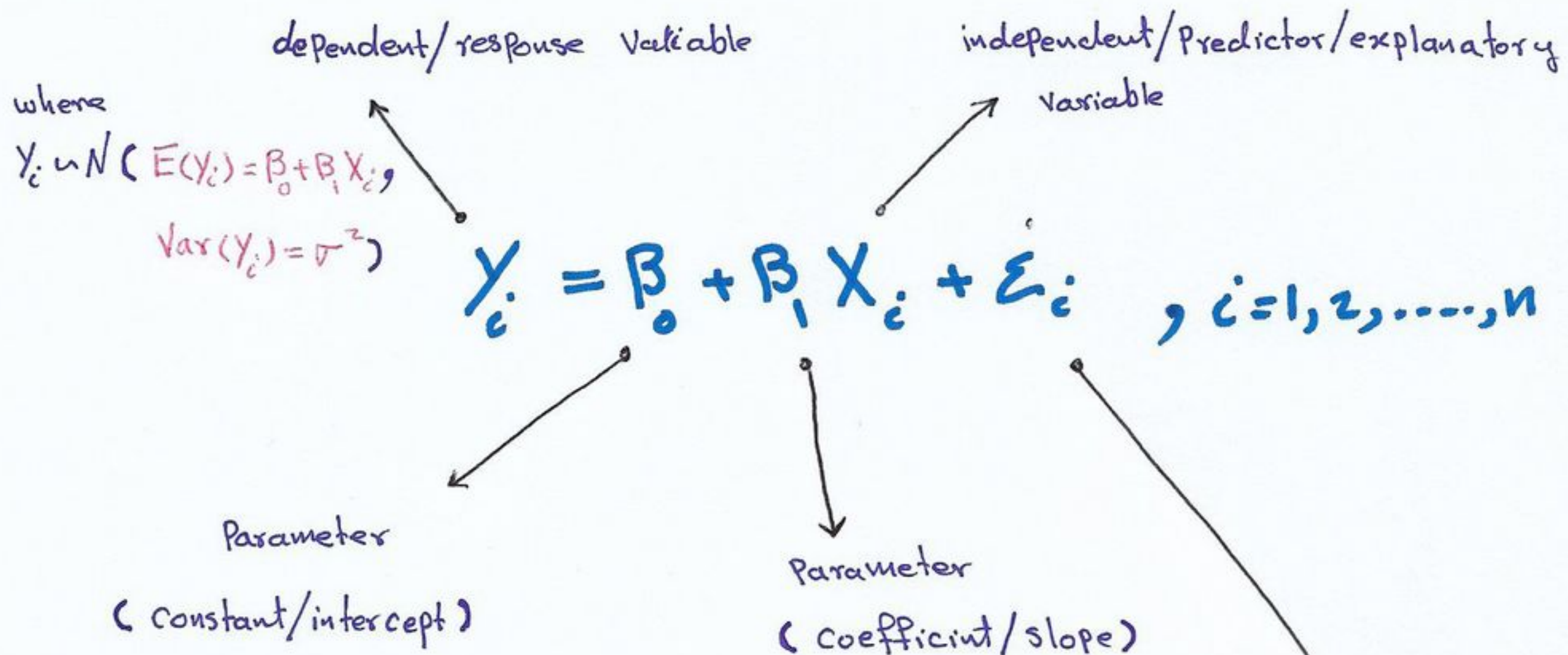
③ (i) Sign of $r$ $\longrightarrow$ $+$ $\longrightarrow$ Positive direction/Correlation   i.e.

   <u>increase</u>-decrease X $\longleftrightarrow$ <u>increase</u>-decrease Y

   $-$ $\longrightarrow$ negative direction/correlation i.e.

   <u>increase</u>-decrease X $\longleftrightarrow$ decrease-increase Y

   (ii) Value of $r$ $\longrightarrow$ $|r| = 0$ $\longrightarrow$ no correlation/ no linear relationship

   $0 < |r| \leq .25$ $\longrightarrow$ weak

   $.25 < |r| < .75$ $\longrightarrow$ moderate

   $.75 \leq |r|$ $\longrightarrow$ strong

dependent / response variable

independent / Predictor / explanatory variable

where

$Y_i \sim N( E(Y_i) = \beta_0 + \beta_1 X_i,$

$Var(Y_i) = \sigma^2 )$

$$Y_i = \beta_0 + \beta_1 X_i + \varepsilon_i \quad , \quad i = 1, 2, \ldots, n$$

Parameter

( constant / intercept )

Parameter

( coefficient / slope )

error

where

$\varepsilon_i \sim N( E(\varepsilon_i) = 0,$

$Var(\varepsilon_i) = \sigma^2 )$

and

$Cov(\varepsilon_i, \varepsilon_j) = 0, \forall i \neq j$

## "Fitted regression line/equation"
### or
### "Prediction equation"

Fitted Value → 

independent/predictor variable →

$$\hat{Y}_i = b_0 + b_1 X_i$$

← constant/intercept

coefficient/slope →

① $b_0 = \begin{cases} \bar{y} - b\bar{x} \\[2mm] \sum_{i=1}^{n} L_i Y_i \; ; \; L_i = \frac{1}{n} - \frac{(X_i - \bar{X})}{S_{XX}} \end{cases}$

② $E(b_0) = \beta_0$ i.e $b_0$ is unbiased estimate of $\beta_0$ ;

$Var(b_0) = MSE\left(\frac{1}{n} + \frac{(\bar{X})^2}{S_{XX}}\right)$ ;

standard error of $b_0$ is $S.E.(b_0) = \sqrt{Var(b_0)}$

③ The sampling distribution of $b_0$:

$T_0 = \frac{b_0 - \beta_0}{S.E.(b_0)} \sim t_{n-2}$

④ $100(1-\alpha)\%$ confidence interval of $\beta_0$:

$\beta_0 \; \epsilon \; b_0 \pm t_{\frac{\alpha}{2}, n-2} \; S.E.(b_0)$

① $b_1 = \begin{cases} \dfrac{S_{xy}}{S_{xx}} = r\sqrt{\dfrac{S_{yy}}{S_{xx}}} \\[3mm] \sum_{i=1}^{n} k_i Y_i \; ; \; k_i = \dfrac{(X_i - \bar{X})}{S_{xx}} \end{cases}$

② $E(b_1) = \beta_1$ i.e $b_1$ is unbiased estimate of $\beta_1$ ;

$Var(b_1) = \frac{MSE}{S_{xx}}$ ;

standard error of $b_1$ is $S.E.(b_1) = \sqrt{Var(b_1)}$

③ The sampling distribution of $b_1$:

$T_1 = \frac{b_1 - \beta_1}{S.E.(b_1)} \sim t_{n-2}$

④ $100(1-\alpha)\%$ confidence interval of $\beta_1$:

$\beta_1 \; \epsilon \; b_1 \pm t_{\frac{\alpha}{2}, n-2} \; S.E.(b_1)$

✱ $b_0$ and $b_1$ are estimate of $\beta_0$ and $\beta_1$, respectively,

by

using

the least squares method / minimization procedure

# Some notes:

① The error/residual $e_i = Y_i - \hat{Y}_i$, then:

    (i) Sum of squared error:
$$SSE = \sum_{i=1}^{n} e_i^2$$

    (ii) $\sum_{i=1}^{n} e_i = 0$

    (iii) mean Squared error:
$$MSE = \hat{\sigma}^2 = S^2 = \frac{SSE}{n-2} \ ;$$

    and $E(S^2) = \sigma^2$ i.e $S^2$ is unbiased estimate of $\sigma^2$

② $\sum_{i=1}^{n} Y_i = \sum_{i=1}^{n} \hat{Y}_i$

③ The sum of weighted residuals:

    (i) by X : $\sum_{i=1}^{n} X_i e_i = 0$

    (ii) by Y : $\sum_{i=1}^{n} Y_i e_i = 0$

④ regression line through $(\bar{X}, \bar{Y})$

⑤ for $k_i$ : (i) $\sum_{i=1}^{n} k_i = 0$      for $L_i$ : (i) $\sum_{i=1}^{n} L_i = 1$

    (ii) $\sum_{i=1}^{n} k_i X_i = 1$             (ii) $\sum_{i=1}^{n} L_i X_i = 0$

    (iii) $\sum_{i=1}^{n} k_i^2 = \frac{1}{S_{xx}}$          (iii) $\sum_{i=1}^{n} L_i^2 = \frac{1}{n} - \frac{(\bar{X})^2}{S_{xx}}$

⑥    $(SSTOT = S_{yy}) = SSE + SSR$

    total sum of squared error

    sum of squared regression
$$= \sum_{i=1}^{n} (\hat{Y}_i - \bar{Y})^2$$

⑦ The coefficient of determination
$$0 \leq R^2 = \frac{SSR}{SSTOT} = 1 - \frac{SSE}{SSTOT} \leq 1$$

    where:

    (i)   or  → $100R^2$% of total Variation in Y is due to X.

    ↘ (the least squares regression/regression line) explains $100R^2$% of the Variation in Y.

    (ii) The simple correlation coefficient $r = \sqrt{R^2}$ and the sign of it will taken from $b_1$.

⑤